

Reviewers' Comments to Author:

(Original comments in italics; our responses below.)

*Some linguistic issues: the nasty neologisms. A "personography" is a biographical list, a "placeography" is a gazetteer. Or all of these things are just "authority files".*

We have inserted the term "authority lists", and also added a footnote explaining our use of personography and placeography in the context of the MoEML project, and how these terms are distinct from "biographical list" and "gazetteer"..

*And if Star Trek uses the word "diagnostics" in a special sense, maybe you should explain it for the non-trekkie reader.*

We have added a footnote explaining how the term is used in Star Trek.

*The term "TEI compliant" doesn't appear anywhere in the literature I'm aware of: the usual term is "TEI conformant".*

The reviewer is correct that "conformance" is more commonly used and is the term used by Burnard 2017. We have corrected throughout.

*The term "referential integrity": although it is used by the MLA, I suspect they appropriated it from database theoreticians, notably C J Date, who wrote several very influential papers on the topic of maintaining the integrity of relational databases.*

Good point. We have updated the footnote. We don't believe it's necessary to research the actual origin of the term, so we have not done so.

*p2 lines 11-12: plural verb singular noun;*

We agree--fixed.

*more to the point your processes \*depend\* on the use of relaxng and schematron don't they?*

We assume that any project is already making full use of regular schema constraints for their XML. We've rephrased to make this clear.

*p4 l10 : I expected to hear more about your tools here, in particular in what way they are "generic" : the problems you're addressing are VERY widespread, and different projects presumably adopt different solutions.*

Our tools are meant to serve as a generic foundation on which projects will build their own set of diagnostic checks. We have now clarified the purpose and the scope of the tools and outlined that projects will need to customize the diagnostics suite to suit their needs.

*Looking at your "referential integrity check" example, and your statement that it's an error if the @type of <name> doesn't match the gi of the target element, I wonder whether that's how you address the issue. Of course if you used <orgName> rather than <name type="org"> the third of your checks would be redundant. I also wonder whether you might wish to say something about the handling of ambiguous cases, e.g. the classic "Elizabeth was very fond of Essex".*

This is an encoding issue, not a diagnostics issue. The specific ambiguity referred to by the reviewer does not occur to our knowledge in any of our texts. The decision to use <name type="x"> rather than the "syntactic sugar" equivalents such as <orgName> is a purely pragmatic one and has no impact on integrity checking; the same check would have to be performed in order to determine whether the encoder was right to use <orgName> versus <persName>. The point is that the encoder makes a claim about what kind of entity they're tagging in their choice of tags, and the diagnostics procedure checks whether the claim is correct.

*Your discussion of level 1 contains a lot of schematron and a discouraging paragraph saying that this solution doesn't scale well. Many readers won't grock the former and will expect some alternative suggestions for the latter. After all there are plenty of projects which get by very nicely whether by using local copies of an authority file, or by using an external database. We agree with the reviewer that not all readers may be familiar with Schematron XML; we have now provided a brief summary of the schematron code block. We also clarified that diagnostics are meant to be supplementary to schematron; for small projects, many checks performed by diagnostics can be done with schematron, but for large projects, diagnostics provide a more efficient and less resource intensive form of validation and consistency checking.*

*p6/l24 "provide progress analysis" noun collision. How about "monitor progress"*  
We agree--changed.

*p7 example : might perhaps be nice to include an example of a real duplicate as well as these faux amis*

*We agree. New image included.*

*p8 "Agas" ???*

We have now given a brief explanation of the column headings, and link to MoEML's explanation of the Agas map.

*pp 8-9 These level 3 features are interesting: their potential use for predicting project end dates seems a little implausible though, especially in view of your subsequent discussion: the better your diagnostics, the more things you will find to diagnose. You should perhaps give some hints as to how to decide whether your current set of diagnostic procedures is adequate.*

*We believe this is addressed in the section "Too much of a good thing".*

-

*(This whole field is one that NLP people worry about endlessly: there are numerous theoretical papers on validation of linguistic corpora.)*

This seems outside the scope of our paper.

*And (contra Kirschenbaum) L Burnard has pointed out several times that "nothing in digital form is ever really finished."*

[We agree that there is a common tacit assumption that digital project are never finished, but we](#)

believe that this is rather a pernicious idea; we have added the Burnard quote to the last paragraph along with some explanation of our position.